# Introduction

The objective of this project is to develop a sentiment analysis model to predict stock price movements based on textual data from news articles, social media posts, and other sources of financial news and opinions. By analyzing the sentiment expressed in these texts, we aim to uncover insights into investor and market sentiment, which can serve as valuable indicators for making informed trading decisions.

# Objectives

1. **Data Collection**: Gather a large dataset of textual data related to stocks, including news articles, social media posts, earnings reports, and analyst reports.
2. **Data Preprocessing**: Clean and preprocess the textual data by removing noise, tokenizing the text into words or phrases, and applying techniques such as stemming and lemmatization.
3. **Labeling Data**: Assign corresponding stock price movements (e.g., increase, decrease, or no change) to the textual data to create a labeled dataset for supervised learning.
4. **Feature Extraction**: Extract features from the textual data, such as word frequencies, sentiment scores, and topic modeling representations.
5. **Model Training and Evaluation**: Train and evaluate various machine learning models, including logistic regression, support vector machines (SVM), random forests, and neural networks.
6. **Performance Metrics**: Use accuracy, precision, recall, F1-score, and receiver operating characteristic (ROC) curve analysis to evaluate model performance.
7. **Trading Strategy Evaluation**: Evaluate the trading strategy using financial metrics such as the Sharpe ratio, maximum drawdowns, number of trades executed, and win ratio.

# Data Collection and Preprocessing

## Data Collection

The dataset comprises textual data related to stocks from various sources, such as:

- News articles from financial news websites
- Social media posts (e.g., tweets mentioning stock tickers)
- Earnings reports
- Analyst reports

## Data Preprocessing

The preprocessing steps include:

1. **Noise Removal**: Removing irrelevant characters, punctuation, and numbers.
2. **Tokenization**: Splitting the text into individual words or phrases.
3. **Stemming and Lemmatization**: Reducing words to their base or root forms to standardize text representations.

# Labeling Data

The textual data is labeled based on the corresponding stock price movements over a specified time horizon. For instance, if a news article is published at time $t$, the stock price movement is observed over the next day, week, or month to determine the label (increase, decrease, or no change).

# Feature Extraction

## Word Frequencies

Term Frequency-Inverse Document Frequency (TF-IDF) is used to represent the text data in numerical form based on word frequencies.

## Sentiment Scores

Sentiment analysis tools (e.g., VADER, TextBlob) are used to calculate sentiment scores for each text entry, indicating positive, negative, or neutral sentiment.

## Topic Modeling

Latent Dirichlet Allocation (LDA) is used to identify topics within the text data, providing additional features for the machine learning models.

# Model Training and Evaluation

## Machine Learning Models

We train several machine learning models, including:

- Logistic Regression
- Support Vector Machines (SVM)
- Random Forests
- Neural Networks

## Performance Metrics

The models are evaluated using the following metrics:

- **Confusion Matrix**: To visualize true positives, false positives, true negatives, and false negatives.
- **Accuracy**: The ratio of correctly predicted instances to the total instances.
- **Precision**: The ratio of correctly predicted positive observations to the total predicted positives.
- **Recall**: The ratio of correctly predicted positive observations to all the observations in the actual class.
- **F1-Score**: The weighted average of Precision and Recall.

- **ROC Curve**: A graphical representation of the true positive rate versus the false positive rate.

## Results

The confusion matrix and performance metrics for the best-performing model are as follows:

- Confusion Matrix:

```lua
Copy code
[[139  47]
 [ 13 179]]
```

- Accuracy: 0.841
- Precision, Recall, and F1-Score:

```markdown
Copy code
              precision    recall  f1-score   support
         0        0.91      0.75      0.82       186
         1        0.79      0.93      0.86       192
micro avg        0.84      0.84      0.84       378
macro avg        0.85      0.84      0.84       378
weighted avg     0.85      0.84      0.84       378
```

# Trading Strategy Evaluation

To evaluate the practical applicability of our model, we simulate a trading strategy based on the predictions and measure the following financial metrics:

## Sharpe Ratio

The Sharpe ratio measures the performance of an investment compared to a risk-free asset, after adjusting for its risk. It is defined as:

Sharpe Ratio=Average Return−Risk-Free RateStandard Deviation of Return\text{Sharpe Ratio}

## Maximum Drawdowns

Maximum drawdown is the maximum observed loss from a peak to a trough of a portfolio, before a new peak is attained. It provides insight into the risk of a trading strategy.

## Number of Trades Executed

The total number of trades executed based on the model's predictions is tracked to evaluate the trading frequency.

## Win Ratio

The win ratio is the proportion of profitable trades to the total number of trades executed.

## Conclusion

The sentiment analysis model demonstrates strong predictive performance with an accuracy of 84.1%, a precision of 0.85, a recall of 0.84, and an F1-score of 0.84. These results suggest that sentiment analysis can be a valuable tool for predicting stock price movements. Further evaluation using financial metrics such as the Sharpe ratio, maximum drawdowns, number of trades executed, and win ratio will provide additional insights into the practical applicability of the model.